1   The evolution of genomic islands by increased establishment probability of
2   linked alleles
3
4   Sam Yeaman [*†‡•], Simon Aeschbacher[§•], and Reinhard Bürger [§§]

5   [†]Biological Sciences, University of Calgary, Calgary, AB T2N 1N4, Canada
6   [‡]Biodiversity Research Centre, University of British Columbia, Vancouver, BC V6T
7   1Z4, Canada [§]Department of Evolution and Ecology, University of California, Davis, CA
8   95616, USA [§§]Faculty of Mathematics, University of Vienna, Oskar–Morgenstern–Platz
9   1, A–1090 Vienna, Austria

10  *samuel.yeaman@ucalgary.ca
11  •*These authors contributed equally*
12
13
14  **Key words:** gene flow, local adaptation, linkage disequilibrium, recombination, genetic
15  architecture, divergence hitchhiking
16
17
18  **Abstract**
19  Genomic islands are clusters of loci with elevated divergence that are commonly found in

20  population genomic studies of local adaptation and speciation. One explanation for their

21  evolution is that linkage between selected alleles confers a benefit, which increases the

22  establishment probability of new mutations that are linked to existing locally adapted

23  polymorphisms. Previous theory suggested there is only limited potential for the

24  evolution of islands via this mechanism, but involved some simplifying assumptions that

25  may limit the accuracy of this inference. Here we extend previous analytical approaches

26  to study the effect of linkage on the establishment probability of new mutations, and

27  identify parameter regimes that are most likely to lead to evolution of islands via this

28  mechanism. We show how the interplay between migration and selection affects the

29  establishment probability of linked vs. unlinked alleles, the expected maximum size of

30  genomic islands, and the expected time required for their evolution. Our results agree

31    with previous studies, suggesting that this mechanism alone is unlikely to be a general

32    explanation for the evolution of genomic islands. However, this mechanism could occur

33    more readily if there were other pre-adaptations to reduce local rates of recombination or

34    increase the local density of mutational targets within the region of the island. We also

35    show that island formation via erosion following secondary contact is much more rapid

36    than island formation from *de novo* mutations, suggesting that this mechanism may be

37    more likely.

38

39    **Introduction**

40    Studies of natural populations commonly find that genetic divergence is elevated in a

41    number of genomically restricted regions along the chromosomes. These have been

42    dubbed "genomic islands" of speciation (Turner *et al.* 2005), differentiation (Harr 2006),

43    or divergence (Nosil *et al.* 2009), depending on the biological context and the way they

44    are measured. Yet, other studies have not found such islands, although they might have

45    been expected in their respective contexts (Nosil *et al.* 2009; Strasburg *et al.* 2012;

46    Renaut *et al.* 2013). What can be learned about the evolutionary history of a species by

47    studying these genomic islands, and what does their absence tell us? There are many

48    different ways that the interplay between selection and demography can give rise to this

49    characteristic pattern in the genome. Therefore, identifying ways to discount alternative

50    explanations and evaluate the significance of presence, absence, and extent of genomic

51    islands is critical to making inferences about evolutionary history.

52        Mechanistic explanations for the evolution of genomic islands can be categorized

53    by whether clusters arise as: 1) purely the result of drift and demography, 2) neutral by-

54    products of linkage and natural selection with or without on-going gene flow, 3) a direct

55    adaptive response, due to the advantage of physical linkage among loci experiencing

56    divergent natural selection and gene flow (summarised in Table 1), or through some

57    combination of these explanations. These three broad categories reflect differences in the

58    number of selected loci in the island that affect the trait under divergent selection, with

59    either none, one, or several causal loci per island, respectively.  We emphasise that

60    explanation 2 can lead to situations where multiple selected loci, each with their own

61    cluster of neutral sites, are located in close proximity along the genome by random

62    chance. This would lead to an empirical pattern very difficult to distinguish from

63    explanation 3 in practice. For conceptual clarity, however, we categorize these

64    mechanisms in terms of the number of causal loci and the type of linkage-mediated

65    benefit among them.

66         Before reviewing these three mechanisms, it is important to situate our approach

67    in the context of the many ways that genetic divergence can be quantified and used to

68    identify genomic islands. The most widely used measures of divergence are $F_{ST}$ (Wright

69    1931, 1943) and related statistics (Nei 1973; 1982; Slatkin 1991; Charlesworth 1998;

70    Excoffier 2007), which express the variance in allele frequency among (partially) isolated

71    populations, normalized by the genetic diversity in an appropriately defined panmictic

72    population. As such, $F_{ST}$-like statistics are relative measures of divergence. They are

73    sensitive to a variety of processes affecting diversity within and among populations, and

74    therefore tend to be easily confounded. In contrast, absolute measures of divergence such

75    as $D_{xy}$ (Nei & Li 1979) are more specific, but not very sensitive with certain types of

76    data. Here, we formulate our discussion of islands of divergence in terms of $F_{ST}$. Where

77    necessary, we mention potential confounding factors, and in Table 1 we contrast $F_{ST}$ to

78    patterns including absolute measures that are diagnostic for various mechanisms. For a

79    review of relative and absolute measures and their properties, we refer to Cruickshank &

80    Hahn (2014).

81        Although the action of natural selection through either mechanism 2 or 3 may feel

82    intuitively most parsimonious, it is important to discount neutral explanations in

83    empirical studies (i.e. mechanism 1, Table 1). If a pattern of genomic islands is found

84    along a single gradient or between a single pair of populations, it may be difficult to

85    conclusively rule out the effect of demography. The distribution of coalescence times for

86    neutral alleles becomes much broader under isolation-by-distance or hierarchical

87    structure than under panmixia, so loci that appear to be outliers under an island model

88    may in fact still be consistent with drift (Excoffier *et al.* 2009; Hermisson 2009). Linkage

89    among loci in close proximity on a chromosome could then cause the appearance of an

90    island of extreme $F_{ST}$. Additionally, rapid population expansion can result in 'allelic

91    surfing' due to increased genetic drift at the wave front (e.g. Slatkin & Excoffier 2012).

92    This could likely cause blocks of elevated divergence, especially in regions of low

93    recombination, but has not been extensively studied. However, if the same genomic

94    islands are found along several demographically independent gradients, it is much less

95    likely that the pattern could have arisen purely by drift and demography.

96        Within the second category of explanation, there are several different ways that

97    linkage of neutral loci to an allele under selection can result in elevated divergence along

98    a chromosome. Recurrent bouts of either purifying or positive and spatially homogenous

99    selection can result in reduced genetic variation within populations at regions linked to

100 the loci under selection (Maynard Smith & Haigh 1974; Charlesworth *et al.* 1993), which

101 can inflate relative metrics of divergence, such as $F_{ST}$ (sub-category 2A in Table 1; Nei

102 1973; Charlesworth *et al.* 1997; Nordborg 1997; Cruickshank & Hahn 2014). It is unclear

103 whether alternative, absolute metrics of divergence that are insensitive to variation within

104 populations, such as $D_{xy}$, are sensitive enough to address this problem, but if genomic

105 islands do not co-occur with locally reduced heterozygosity, then the explanations of

106 purifying selection and global sweeps can be discounted. Alternatively, spatially

107 heterogeneous selection can lead to increased divergence at linked neutral regions. If this

108 occurs as a single selective sweep within a restricted region of the species range (without

109 gene flow), then hitchhiking can decrease heterozygosity and increase divergence at

110 linked sites, but this signature will decay over time (sub-category 2B in Table 1; Maynard

111 Smith and Haigh 1974; Gillespie 2000). When divergent selection operates in the face of

112 sufficiently low gene flow, such local sweeps are never completed and an equilibrium is

113 reached where the change in allele frequency due to selection is opposed by gene flow. In

114 populations of finite size, this recurrent selection against maladaptive gene flow generates

115 a persistent pattern of increased divergence and reduced heterozygosity at linked sites

116 (sub-category 2C in Table 1; Charlesworth *et al.* 1997; Nordborg 1997; Via & West

117 2008; Feder & Nosil 2010). Finally, if locally adapted and previously allopatric

118 populations come into secondary contact, erosion of divergence will occur most rapidly at

119 neutral loci that are less tightly linked to a given adapted locus (sub-category 2D in Table

120 1; Petry 1983; Barton and Bengtsson 1986; Barton and Hewitt 1989).

121     By the third category of explanation, genomic islands evolve as an adaptive

122 consequence of the way recombination mediates the tension between divergent selection

123  and gene flow. If two or more locally beneficial mutations are tightly physically linked,

124  recombination is less likely to break up this favourable combination than when linkage is

125  weak (Fisher 1930; Charlesworth 1979; Lenormand & Otto 2000). This advantage due to

126  linkage can lead to clustered genetic architectures under divergent selection in the face of

127  gene flow, through one of four potential evolutionary routes (category 3 in Table 1): A)

128  increased probability of establishment of linked alleles during initial divergence (Yeaman

129  & Whitlock 2011; Feder *et al.* 2012; Aeschbacher & Bürger 2014); B) competition

130  among genetic architectures as new linked mutations displace older unlinked mutations of

131  equal overall effect (Yeaman and Whitlock 2011); C) competition among genomic

132  architectures, which favours the fixation of rearrangements moving loci into close

133  proximity on a chromosome (Yeaman 2013) or the establishment of segregating

134  inversions that reduce recombination (Noor *et al.* 2001; Rieseberg 2001; Kirkpatrick

135  2006); D) increased persistence time for linked and selected loci under secondary contact

136  following local adaptation in allopatry (Petry 1983; Barton & Bengtsson 1986; Barton &

137  Hewitt 1989). In all of the above cases in category 3, selection at the causal loci would

138  also cause elevated divergence at linked neutral loci.

139       One way to differentiate between explanations of category 2 and 3 is to identify

140  whether, in a single genomic island, there are multiple sites harbouring alleles that

141  contribute to fitness differences between populations. This is expected to occur under the

142  latter but not the former explanation. While it is very difficult to identify which mutations

143  are functionally important, it may be possible to infer their presence from observations

144  about multiple peaks within an island, or through inferences about the maximum

145  expected width of an island under explanations of category 2 vs. 3. Furthermore, we

146    could learn much about a species' evolutionary history by discriminating between the

147    four possible explanations within category 3. In practice, this is difficult because some of

148    the expected statistical signatures from these routes to clustering may be very similar at

149    equilibrium. For instance, at equilibrium, both scenarios 3.A) and 3.D) lead to an

150    expected pattern of $F_{ST}$ peaks surrounded by regions of low $F_{ST}$. That said, the

151    trajectories to this equilibrium pattern may be quite different. In 3.A), peaks of

152    divergence could be seen as mountains arising against a more or less constant

153    background, whereas in 3.D) regions of reduced divergence are better viewed as valleys

154    eroding between mountains of a (roughly) fixed height. As such, it is helpful to use

155    theoretical methods to investigate the expected properties of these various explanations

156    and examine their relative likelihood.

157         Theoretical arguments have been made suggesting that the higher establishment

158    probability (EP) of linked mutations (explanation 3.A) is unlikely to result in a strong

159    signal of clustering, because the increase in EP is small and extends over a small region

160    of a given chromosome, relative to the size of the unlinked portion of the genome

161    (Yeaman 2013). However, the analysis by Yeaman (2013) discounted the effect of

162    linkage in the regions of parameter space where it is potentially most important, so the

163    inferences made about the likelihood of genomic island evolution via increased EP may

164    be incorrect. Moreover, Yeaman (2013) employed a semi-heuristic approximation to the

165    EP (Yeaman & Otto 2011) that has not been formally tested in the case of multiple linked

166    loci under selection. Here, we first verify this approximation by comparison to an

167    alternative, more formally derived one (Aeschbacher & Bürger 2014). Second, we extend

168    the work of Yeaman (2013) to quantify more comprehensively the potential importance

169     of increased EP at linked mutations as an explanation for the evolution of genomic

170     islands. Specifically, we explore predictions about the expected width of a genomic

171     island, the absolute and relative increase in EP due to linkage between selected sites, the

172     waiting time for island evolution, and the effect of different distributions of mutational

173     effect size on these predictions. Third, we determine how fast neutral divergence that has

174     built up between two loci under spatially divergent selection without gene flow erodes

175     upon secondary contact, and compare this to the time it takes for an equivalent two-peak

176     island to arise by divergent selection with on-going gene flow. Overall, our aim is to

177     refine the theoretical arguments about island evolution, provide some testable predictions

178     about island size, and identify how further study could aim to test the other hypotheses

179     for island evolution, particularly through secondary contact and erosion of genetic

180     divergence around selected loci.

181

182

183     **Model**

184     To explore the likelihood of genomic islands evolving by increased probability of

185     establishment of new linked mutations (explanation 3.A, Table 1), it is necessary to

186     consider both the effect of linkage on the establishment probability (EP) and the rate of

187     occurrence of linked vs. unlinked mutations. To study the effect of linkage on the EP of a

188     *de novo* mutation that occurs near an already established locally adapted allele, we

189     describe a model that incorporates the relevant evolutionary forces (*i.e.* selection, gene

190     flow, recombination, and genetic drift), but is still simple enough to allow for efficient

191     mathematical analysis. We then use simplifying assumptions about the expected

192    distribution of fitness effects of new mutations and the relative size of the mutational

193    target in linked vs. unlinked regions to parameterize our model and make inferences

194    about the likelihood of genomic island evolution via increased EP.

195         We consider a discrete-time, monoecious diploid model with continent–island

196    type migration (Haldane 1930; Wright 1931) and divergent selection at two linked

197    biallelic loci A and B. Let the alleles at these two loci be $A_1$ and $A_2$, and $B_1$ and $B_2$,

198    respectively. We define the migration rate $m$ as the proportion of the island population

199    replaced by immigrants from the continent each generation, and denote by $r$ the

200    probability of recombination between the two loci per generation. Assuming additive

201    interactions across alleles and loci, and ignoring epistasis as well as parental and position

202    effects in heterozygotes, we write the fitnesses of the nine distinguishable genotypes in

203    the island population as

204
$$\begin{array}{c} \begin{array}{ccc} B_1 B_1 & B_1 B_2 & B_2 B_2 \end{array} \\ \begin{array}{c} A_1 A_1 \\ A_1 A_2 \\ A_2 A_2 \end{array} \begin{pmatrix} 1+a+b & 1+a & 1+a-b \\ 1+b & 1 & 1-b \\ 1-a+b & 1-a & 1-a-b \end{pmatrix} \end{array},$$
(1)

205

206    where $a$ and $b$ are the selective advantages on the island of alleles $A_1$ and $B_1$ relative to $A_2$

207    and $B_2$, respectively. To enforce positive fitnesses, we require that $0 < a, b < 1$, and $a + b$

208    $< 1$. Unless otherwise stated, we assume that the continental population is fixed for alleles

209    $A_2$ and $B_2$, whereas alleles $A_1$ and $B_1$ are local *de novo* mutations endemic to the island

210    population. The conditions for establishment and maintenance of one- and two-locus

211    polymorphisms in the island population have been described before (Bürger & Akerman

212    2011; Aeschbacher & Bürger 2014). For mathematical convenience, we assume that

213  selection in favour of $A_1$ is weaker than selection in favour of $B_1$ ($a < b$). This assumption

214  implies that allele $A_1$ cannot be maintained in the island population if allele $B_1$ is

215  swamped by gene flow (Haldane 1930; Bürger & Akerman 2011). It is further in

216  agreement with our principal focus on the effect of linkage on establishment of a weak,

217  locally beneficial *de novo* mutation arising in the presence of an already existing

218  migration–selection polymorphism at a background locus. Throughout, we assume that

219  populations are large enough such that genetic drift can be ignored after an initial

220  stochastic phase during which *de novo* mutations arise in low absolute numbers. In this

221  setting, locally beneficial *de novo* mutations have a strictly positive establishment

222  probability if they can invade in an equivalent fully deterministic model, and *vice versa*

223  (Bürger & Akerman 2011; Aeschbacher & Bürger 2014). Therefore, if in the following

224  we say that "$A_1$ can be established", this means both that $A_1$ has a strictly positive

225  establishment probability and that it will invade in a corresponding deterministic model.


226       Aeschbacher & Bürger (2014) derived several results for this model that will help

227  interpret our findings discussed below. We therefore briefly recapitulate them in the

228  following. First, in a single-locus model (no background locus B), $A_1$ can be established if

229  and only if $m < a$ (cf. Haldane 1930). Second, if $m \geq a$, the additional effect of linkage to

230  a background locus B means that $A_1$ can be established if $m < m^*$, where $m^* =$

231  $\frac{a(b-a+r)}{(a-r)(a-b)+r(1-a)}$ (their Eq. 10). Third, this critical threshold can alternatively be expressed

232  in terms of a critical recombination rate, where linkage must be at least as tight as $r^* =$

233  $\frac{a(a-b)(1+m)}{a(1+2m)-(1+b)m}$ when $m > a / (1 - 2a + b)$ (their Eq. 11). Fourth, if $m > b/(1 - a)$, then no

234  polymorphism will be maintained at locus B, which automatically implies that $A_1$ cannot

235    be established because gene flow is too strong.

236        Based on this general model, we derive three approximations for the expected EP

237    of any given new mutation, $\pi = \pi(a, r)$, expressed as a function of $a$ and $r$, which are

238    discussed in more detail below and in the supplementary materials. To make predictions

239    about the average establishment probability over all available linked mutations, $\bar{\pi}_L$, we

240    integrate $\pi(a, r)$ across the range of possible recombinational distances $r$ at which new

241    mutations could be linked to the background locus, and across the expected Distribution

242    of Fitness Effects (DFE) for values of $a$. We assume that the DFE is a gamma

243    distribution,

244        $$f_a(k, \bar{a}) = \frac{\left(\frac{k}{\bar{a}}\right)^k}{\Gamma(k)} a^{k-1} e^{-\left(\frac{k}{\bar{a}}\right)a},$$        (2)

245    with shape parameter $k$ and mean $\bar{a}$, where $k, \bar{a} > 0$. If $k = 1$, this simplifies to an

246    exponential distribution, $f_a(1, \bar{a}) = \frac{1}{\bar{a}} e^{-a/\bar{a}}$. With $a$ drawn according to equation (2), our

247    assumption of $a < b$ does not necessarily hold; in practice, we therefore set $\bar{a} \ll b$. When

248    integrating over $r$, we assume a uniform distribution of $r$ between 0 and ½, usually using

249    an upper limit of $r_f$ = ½, because this represents the threshold to free recombination in a

250    discrete-time model. This approach to integration is biologically akin to assuming that

251    mutations occur with equal probability at all positions within the window of the

252    chromosome where $0 < r < r_f$, and the rate of recombination does not vary along the

253    chromosome. For mutations with $r = 1/2$ on the same chromosome, they are essentially

254    unlinked and are not assumed to contribute to $\bar{\pi}_L$. To represent the average EP of an

255    unlinked mutation, $\bar{\pi}_U$, we multiply $\pi(a, r)$ by the DFE, and integrate over $a$, setting $r = r_f$

256    = ½. To obtain the approximate size of a genomic island evolving via the benefit of

257    selection at a linked background locus (mechanism 3.A, Table 1), we determine the

258    region of the chromosome over which most new mutations will experience an increased

259    EP via this mechanism.  Specifically, we calculate the value of $r$ required to contain 95%

260    of the probability density of $\pi(a, r)$, and refer to this as the 95% window size, or $C_{95}$.

261        To explore the time scales over which the last step in the evolution of two-peak

262    genomic islands occur under alternative scenarios (explanations 3.A vs. 3.D), we use a

263    combination of stochastic and deterministic theory. We extend our model to bidirectional

264    gene flow and predict the dynamics of neutral divergence using the structured coalescent.

265    Thereby, we replace the neutral migration rate by appropriate rates of effective gene flow,

266    which are a function of the location of the neutral site and the strength of selection at the

267    two loci causing the adaptive peaks (see SI for details).

268

269    **Results**

270    *Analytical predictions for establishment probability*

271

272    *Two-type branching process approximation in discrete time*: The EP under a two-type

273    branching process, $\pi(a, r)$, was derived in Aeschbacher & Bürger (2014), and represents

274    a weighted average over $A_1$ occurring either on the locally beneficial ($B_1$) or deleterious

275    ($B_2$) genetic background (see SI for details). The EP, $\pi(a, r)$, decreases both with the

276    migration rate $m$ and the recombination $r$, and increases with the selection coefficient $a$

277    (solid lines in Figure 1).  As expected, if migration is sufficiently weak, then $A_1$ can

278    establish even if unlinked to the background locus B (solid blue line in Figure 1B).

279     To represent the average probability of a new, linked mutation across all possible

280     recombination rates and sizes of selection coefficient, we integrate across both the rate of

281     recombination and the DFE,

282

283     $\bar{\pi}_L = \int_0^\infty \int_0^{r_f} \pi(a, r) f_a(k, \bar{a}) f_r(0, r_f) dr\, da,$           (3)

284

285     where $f_r(r_{min}, r_{max})$ is a uniform density between $r_{min}$ and $r_{max}$. Other distributions

286     could be used to represent non-homogenous recombination rates or mutation targets, but

287     for simplicity we focus on the uniform case here. We could not find a simple closed form

288     for $\bar{\pi}_L$, so in the following analyses we use numerical integration to investigate how it

289     varies with $\bar{a}$, $b$, $r$, $m$, and $k$. The difference between the EP of an average linked vs. an

290     average unlinked mutation of a certain effect size $a$, i.e. $\bar{\pi}_L^{(r)}(a)$ vs. $\bar{\pi}_U^{(r)}(a)$, increases with

291     $m$, and above the critical threshold $m^*$, the EP for an individual unlinked mutation is 0

292     (Figure 2A). Integrating across the DFE, and hence accounting for all possible values of

293     $a$, shows that the mean EP of unlinked mutations, $\bar{\pi}_U$, decreases more gradually and does

294     not rapidly intersect 0, as it includes the effects of rare large mutations that have non-zero

295     EP (Figure 2B).

296

297     *Slightly-supercritical branching process*: To derive an approximation to the branching-

298     process solution described above, we assume that the selective advantage of the focal

299     mutation and the migration rate are both small relative to the selection coefficient at the

300     background locus and the recombination rate, *i.e.* that $a$, $m \ll b$, $r$. Under these

301     assumptions, it can be shown that the EP of a linked mutation is

302

303 $$\pi(a, r) \approx 2 \frac{a(b+r) - (1+b)mr}{(1+b)(b+r)}$$ (4)

304 if $m < a(b + r) / [r(1 + b)]$, and 0 otherwise (see SI for derivation). By this

305 approximation, mutations can establish if the recombination rate $r$ is below a critical

306 value given by

307 $$\tilde{r} = \frac{ab}{(1+b)m - a}.$$ (5)

308 If we further assume that both selection coefficients are much smaller than unity, *i.e.* $0 <$

309 $a \ll b \ll 1$, then this reduces to $\tilde{r} = ab/m$, which is identical to the approximate critical

310 recombination rate for invasion of $A_1$ in a deterministic continuous-time model (see Eq.

311 4.15 in Bürger & Akerman, 2011). This can be further reduced to $\tilde{r} \approx b$ when $m \sim a$,

312 which is Barton's rule of thumb that selection against gene flow has an appreciable effect

313 on linked sites only if linkage is sufficiently strong, i.e. $r \ll b$ (Barton 2000).

314 Under the assumption of an exponential DFE ($k = 1$), a closed-form solution for the

315 average EP with respect to $r$ and $a$ can be derived by using a Taylor series approximation

316 around $b = 0$ (see SI for details). This yields

317 $$\bar{\pi}_L \approx e^{-\frac{m}{a}} \left\{ \bar{a} + \bar{a}b + bm - 2bm\gamma + 2bm \left[ E_i\left(\frac{m}{a}\right) + \ln\left(\frac{\bar{a}}{2bm}\right) \right] \right\} - 2\bar{a}b,$$ (6)

318 where $\gamma = 0.577$ is the Euler–Mascheroni constant and $E_i(z) = -\int_{-z}^{\infty} \frac{e^{-t}}{t} \, dt$ is the

319 exponential integral function. This approximation works well as long as $b \lesssim 0.2$ (see SI

320 and Figure S1 for details), and is very close to the splicing approximation (see below).

321    Both of these approximations overestimate the exact two-type branching process when

322    either *m* is large or *r* is small (Figure 1).

323

324    *Splicing approximation in discrete time*: The rate of increase in frequency of allele $A_1$

325    when rare is given approximately by the leading eigenvalue of the Jacobian matrix

326    describing the stability around the equilibrium where B is polymorphic and A is fixed for

327    $A_2$ (see SI for details; as per Yeaman and Otto 2011). Because this deterministic rate of

328    increase in frequency is equivalent to the effect of natural selection in a single-locus one-

329    population model, this rate can be spliced into the standard probability of fixation for a

330    new mutation ($\Pr[\textit{fix}] = 2s$; Haldane 1927), which yields an approximation to the EP of

331    $A_1$,

332    $\pi(a, r) \approx 2 \max[0, \frac{2 + b - r + m(2a - b - r) + \sqrt{R}}{2(1 - a + b)} - 1],$             (7)

333    where

334    $R = (1 + m)\{b^2(1 + m) + 2b(1 - m)r + r[r - m(4 - 4a - r)]\}.$ (8)

335

336    As in the case of the two-type branching process, we must resort to numerical integration

337    to obtain approximations to the average probability of establishment of linked alleles for

338    a given selection coefficient *a*, $\bar{\pi}_L^{(r)}(a)$, and across the entire DFE, $\bar{\pi}_L$. As shown in

339    Figure 1, the splicing approach (dotted lines) yields an approximation to the EP very

340    close to the one obtained from the slightly-supercritical branching process assuming

341    $a, m \ll b, r$ (dashed lines; Eq. 4). Accordingly, it deviates most from the exact two-type

342    branching process (solid lines; Eq. S3) if *a* is relatively large and, at the same time, either

343     *m* is large or *r* is small. In these cases, the contribution of mutations that occur on the

344     deleterious genetic background $B_2$ becomes important; both the splicing approach and the

345     slightly-supercritical branching process do not capture this.

346

347     *Island size and establishment probability of linked vs. unlinked mutations*

348     We now use the above approximations to explore the likelihood of genomic islands

349     evolving by mechanism 3.A (Table 1). As outlined above, we integrate both across all

350     possible recombination rates and mutation effect sizes (assuming a gamma-distributed

351     DFE with a mean effect size of $\bar{a}$, as above) to represent how an average linked vs.

352     unlinked mutation would be affected by the interplay between migration, selection, and

353     recombination. With increasing migration rate in the region of $m \sim \bar{a}$, an increasing

354     fraction of the available mutations become more strongly limited by migration than

355     favoured by selection. This causes a decrease in the mean EP of both linked ($\bar{\pi}_L$) and

356     unlinked ($\bar{\pi}_U$) mutations, but the effect of linkage to the **B** locus mitigates this effect, so

357     that the mean EP decreases faster for unlinked mutations. These broad patterns are

358     illustrated in Figures 3A & B, which show the effect of *m* on $\bar{\pi}_L$, and Figure 3C, which

359     shows the ratio $\bar{\pi}_L/\bar{\pi}_U$. In each case, there is a pronounced transition in the region of $m \sim$

360     $\bar{a}$. The decrease in $\bar{\pi}_U$ occurs more rapidly than $\bar{\pi}_L$ because as the migration rate

361     increases, $\pi_U$ approaches zero for an increasing fraction of unlinked mutations, whereas

362     linked mutations still have non-zero probability of establishment (cf. Figure 2B).

363     Therefore, $\bar{\pi}_L/\bar{\pi}_U$ increases towards infinity if *m* is increased far beyond $\bar{a}$, and while

364     linked mutations can still establish, their absolute EP is much reduced.

365        Another parallel effect of increasing migration is to reduce the size of the linked

366        region that can contribute to adaptation. As migration increases, an increasing fraction of

367        the chromosome falls above the critical recombination threshold $r^*$, such that an

368        increasing fraction of the available linked mutations cannot be established. As a

369        consequence, there is a pronounced decrease in the expected size of the region around the

370        background locus that can contribute to local adaptation if $m > \bar{a}$, as shown by $C_{95}$, the

371        95% window size (Figure 3D).

372        The shape of the gamma distribution ($k$) affects the position of the transition zone

373        where i) the EP starts decreasing, ii) the ratio $\bar{\pi}_L/\bar{\pi}_U$ starts increasing, and iii) the size of

374        the 95% window starts decreasing. These transitions occur at lower values of $m$ for $k = 2$

375        (*i.e.* a DFE with relatively more intermediate values of $a$, compared to $k = 1$) and higher

376        values of $m$ for $k = 0.5$ (i.e. a DFE with relatively more extreme values of $a$, compared to

377        $k = 1$) (Figure S2). Thus, while $k$ does not affect the broad qualitative patterns

378        representing the role of linkage, it does affect the relative importance of linkage at a

379        given value of $m$ (*i.e.* the broad patterns are maintained, but shifted with respect to $m$ with

380        changes in $k$). Specifically, the average EP of linked mutations, $\bar{\pi}_L$, is higher for smaller $k$

381        because a larger fraction of the probability density occurs with mutations of larger size,

382        which have proportionally larger EP $\bar{\pi}_L^{(r)}$ (Figure S2B). This is also reflected in the width

383        of the genomic window within which most linked mutations establish: smaller $k$ and

384        hence relatively more mutations of large effect increase the size of the window ($C_{95}$ in

385        Figure S2D). However, put in relation to the average EP of unlinked mutations, $\bar{\pi}_U$, the

386        effect of the shape is compensated for by the fact that unlinked mutations profit

387        proportionally more from a fat right-hand tail (low $k$) of the input DFE (Figure S2C).

388    It is also worth noting that the EP is lower for higher values of $b$ when $m$ is very

389    low (see SI for details and a derivation of the critical migration rate). This occurs because

390    the contribution to the total fitness made by locus A becomes relatively smaller as $b$

391    increases (*i.e.*, the marginal fitness ratio of genotype $A_1A_2B.B.$ vs. $A_2A_2B.B.$ decreases

392    with increasing $b$), which reduces the net advantage of allele $A_1$, and hence its EP.

393    However, this effect of the relative magnitude of selection coefficients is only important

394    when the facilitating effect of linkage is relatively small (*a* large relative to $b$); when

395    linkage is important, then the EP increases with $b$.

396

397    *Waiting time for evolution of genomic islands*

398    Taken together, the above results show that linkage becomes critically important for any

399    local adaptation to occur when the migration rate increases beyond the mean selection

400    coefficient ($m > \bar{a}$). However, they also illustrate that in this same region of parameter

401    space, there is a decrease in both the size of the region that can contribute to this

402    adaptation (Figure 3D) as well as the absolute EP of mutations within this region (Figure

403    3A &B). These latter two factors would be expected to greatly increase the waiting time

404    for the first linked mutation to establish, as the mutational target is smaller with small

405    window sizes, and waiting time scales with the inverse of the EP. To obtain a rough

406    approximation to this expected waiting time, we combine the results from Figure 3,

407    calculating $t_L \approx 1/(2N\mu 2C_{95}\bar{\pi}_L)$ and $t_U \approx 1/(2N\mu C_U\bar{\pi}_U)$, for linked and unlinked

408    mutations, where $N$ is the size of the diploid population, $\mu$ is the mutation rate per cM per

409    gamete per generation, $C_{95}$ is the 95% window size in cM as defined above, and $C_U$

410    is the size of the genome that is unlinked to $\mathbf{B}$, in cM. We multiply $C_{95}$ by 2 to account

411    for the fact that a linked mutation can occur on both sides of the background locus.

412    Unfortunately, it is difficult to parameterize these equations to make quantitative

413    predictions about the expected waiting time, because we know little about the rate of

414    beneficial mutations. In organisms such as humans, where genome size is on the order of

415    $10^9$ base pairs (bp), and where there are approximately $10^6$ bp/cM and about $10^{-8}$

416    mutations/bp/generation, the mutation rate per cM is $\sim 0.01$ per generation, but it is

417    unclear what fraction of these would be locally beneficial.

418           Given uncertainty about $N$ and $\mu$, and accounting for appropriate mutational

419    target, we scale the waiting times by the inverse of the substitution rate of an average

420    mutation from the same DFE that arises in a completely isolated panmictic population of

421    size $N$. Because establishment is independent of recombinational distance from the

422    background locus if $m = 0$, this reference mutational target simply amounts to

423    $2C_{95,m=0} = 2 \times 0.95 \times r_f = 0.95$ for $r_f = 0.5$, independently of all the other

424    parameters. Specifically, we approximate the scaled expected waiting times for linked

425    and unlinked mutations as

426

427    $$T_L \approx t_L \left( \frac{1}{4N\mu\bar{a}\,2C_{95,m=0}} \right)^{-1} \approx \frac{4N\mu\bar{a}\,2C_{95,m=0}}{4\,N\mu\,C_{95}\bar{\pi}_L} = \frac{0.95\,\bar{a}}{C_{95}\bar{\pi}_L}, \qquad\qquad (9a)$$

428    $$T_U \approx t_U \left( \frac{1}{4N\mu\bar{a}\,2C_{95,m=0}} \right)^{-1} \approx \frac{4N\mu\bar{a}\,2C_{95,m=0}}{2\,N\mu\,C_U\bar{\pi}_L} = \frac{1.9\,\bar{a}}{C_U\bar{\pi}_U}. \qquad\qquad (9b)$$

429

430    This scaling has the advantage of focussing exclusively on the *relative* impact on waiting

431    times of the tension between migration and divergent selection, and therefore seems

432    appropriate for a comparison of linked vs. unlinked locally beneficial *de novo* mutations.

433         Some uncertainty remains about $C_U$, the size of the unlinked mutational target.

434    However, we can make reasonable choices. For $b = 0.1$ and a $C_U$ of 100 cM, i.e. equal to

435    twice the map distance with respect to which we defined $C_{95}$, the difference in waiting

436    times for linked and unlinked mutations is negligible if $m < \bar{a}$ (dotted vs. solid curves in

437    Figure 4A). This is expected, because in the limit of no migration, the background locus

438    is fixed for $B_1$, and establishment of the focal mutation $A_1$ does not depend on its

439    recombinational distance from locus B. As $m$ increases above $\bar{a}$, $T_U$ increases more

440    rapidly than $T_L$, because linkage now increases the relative establishment probability and

441    the mutational target for linked mutations ($C_{95}$) is still large (compare Figures 3C and

442    3D). At $m \approx 10\bar{a}$, the increase in $T_L$ becomes log-log linear, while $T_U$ keeps increasing

443    exponentially with $m$. This is where linkage starts to convey a considerable relative

444    advantage. However, absolute waiting times for linked mutations are at least about 100

445    times higher compared to the case of $m = 0$ (Figure 4A). This log-log linear phase ends if

446    $m \approx b$, at which point gene flow swamps even the beneficial background allele $B_1$, and

447    $T_L$ quickly increases to infinity.

448         Realistically, the mutational target size for unlinked mutations is much larger than

449    100 cM in most organisms. As expected, for $C_U > 100$ cM, the waiting time to

450    establishment for unlinked mutations becomes much reduced compared to linked

451    mutations if $m$ is weak enough, suggesting that in this region of parameter space, a strong

452    signature of clustering would be unlikely to evolve. The larger the unlinked target size,

453    the higher the migration rate at which linked mutations are expected to arise earlier than

454    unlinked ones (i.e. where dashed curves intersect with solid ones in Figure 4A). When

455    this occurs, $t_L$ tends to be very long (relative to the waiting time in the absence of

456    migration) except at the highest values of $\bar{a}$ (blue curves in Figure 1A). This illustrates

457    that in the region where a strong signal of clustering is expected via mechanism 3.A, the

458    waiting time for the evolution of this pattern will be very long relative to adaptation in

459    allopatry. The shape of the gamma distribution ($k$) has a strong effect on absolute waiting

460    times, but does not substantially affect relative difference between linked and unlinked

461    locally beneficial *de novo* mutations (Figure S3). However, there is some effect of $b$, as

462    the transition to long waiting times with increasing $m$ for linked mutations is sharper and

463    occurs at lower $m$ with small $b$ (Figure 4B).

464

465    *Divergence with gene flow vs. secondary contact*

466    Genomic islands consisting of multiple causal loci can evolve under various mechanisms

467    summarised in category 3 (Table 1). The same amount of divergence is expected at

468    equilibrium in models of secondary contact and *de novo* adaptation (Bürger & Akerman

469    2011; Aeschbacher and  Bürger 2014), which makes it challenging to distinguish between

470    these alternative explanations based on empirical observations. However, by studying the

471    dynamics leading to these equilibrium patterns, it may be possible to exclude some

472    mechanisms based on a comparison of associated time scales with the demographic

473    history of populations and species of interest. To illustrate this argument, we concentrate

474    on the competing mechanisms of increased establishment probability of linked mutations

475    in the face of gene flow (explanation 3.A in Table 1) vs. the erosion of neutral divergence

476    upon secondary contact (explanation 3.D). Moreover, we extend our model to allow for

477    symmetric gene flow between the two demes of equal size, as compared to the

478    asymmetric continent–island regime assumed above. We restrict our treatment to

479    genomic islands made up of two causal loci, and we focus on the last evolutionary step

480    needed to complete equivalent two-locus islands at the equilibrium of migration,

481    selection, and genetic drift. By equivalent we mean that the map distance between the

482    two loci A and B under selection, the migration rate $m$ at parapatric stages, and the

483    selection coefficients $a$ and $b$ are identical in the two scenarios.

484         Under explanation 3.A, we define the last evolutionary step as the rise of a peak

485    of divergence at locus A by establishment of a weak, locally beneficial mutation in one of

486    the two demes. Here, we condition on there being an established migration–selection

487    polymorphism at the linked background locus B, and on the new mutation being

488    successful. We start counting time when the locally beneficial mutation occurs. As above,

489    we focus on the case where linkage to B is essential for establishment of the new

490    mutation. In contrast, under explanation 3.D, we view the erosion of neutral divergence

491    as the last evolutionary step during secondary contact following adaptive divergence at

492    loci A and B in allopatry (assuming the period of allopatry has been long enough to allow

493    local adaptation, which is not included in the following estimates). Here, we condition on

494    locally adaptive mutations having been fixed before secondary contact, and we start

495    counting time at the onset of secondary contact.

496         In both cases, we ignore mutation at loci A and B, and we assume that neutral

497    mutation rates are low, so that the infinite-alleles model of mutation holds and $F_{ST}$ can be

498    approximated as

499

500 $$F_{ST} \approx \frac{T_T - T_W}{T_T} \qquad\qquad\qquad (10)$$

501

502   (Slatkin 1991, Wilkinson-Herbots 1998) Here, $T_T$ and $T_W$ are the expected coalescence

503   times between two samples taken at random from the entire population and from within a

504   deme, respectively. These times depend on the local effective population size $N_e$ and the

505   migration rate $m$ in a neutral model. We derive these coalescence times in the SI using

506   classical results for the structured coalescent (Griffiths 1981; Takahata 1988; Notahara

507   1990; Wilkinson-Herbots 1998, 2008, 2012). To account for the effect of selection, we

508   replace $m$ by the appropriate effective rate of gene flow $m_e$ (Petry 1983; Aeschbacher &

509   Bürger 2014; Akerman & Bürger 2014).

510        Under these assumptions, erosion of neutral divergence during secondary contact

511   after adaptive divergence occurs on a much faster time scale than the rise of a new peak

512   in $F_{ST}$ under *de novo* selection with gene flow (Figure 5). With $a = 0.01$, $b = 0.1$, and $m =$

513   0.02, a potentially identifiable two-peak island is formed already after about $t = 100$

514   generations in the secondary-contact scenario (pale-blue curves in Figure 5A), whereas

515   the peak arising at locus A is barely noticeable at this time; only after about $t = 1000$

516   generations is there a clear second peak (Figure 5B). Under purely deterministic

517   dynamics, i.e. assuming infinitely large local population sizes, the waiting time for the

518   decay of neutral divergence in terms of $F_{ST}$ is well approximated by the inverse of the

519   effective rate of gene flow, $1/m_e$ (SI; Figure S4). Even with a finite local population size

520   $N_e$, $1/m_e$ seems to provide a good approximation to the time until $F_{ST}$ reaches its

521   equilibrium in the secondary-contact case (Figure 5C). In contrast, it takes much more

522    time for $F_{ST}$ at locus A to rise upon occurrence of a locally beneficial mutation in the

523    presence of gene flow (Figure 5D and Figure S5B and S5D).

524        The large differences in time scales over which equivalent genomic islands arise

525    and become potentially detectable suggest that explanation 3.D may be much more likely

526    than 3.A if the two populations or species of interest have diverged relatively recently.

527    Given that the above analysis assumes that a mutation destined to establish has already

528    occurred under 3.A, the total waiting time for mechanism 3.A would be even longer than

529    shown in Figure 5 (as per Figure 4). These differences remain even if adaptive divergence

530    before secondary contact has not been complete, as long as the allopatric phase has not

531    been too short (Figure S5A and S5C). We note that the mutations contributing to local

532    adaptation in allopatry (mechanism 3.D) are not expected to be more clustered than

533    random, unless the underlying mutational target is itself clustered (e.g., due to previous

534    evolution via mechanism 3C). Thus, mechanism 3.D also depends on multiple mutations

535    occurring in close enough proximity to each other that a contiguous island is observed, at

536    least for some time after secondary contact.

537

538    **Discussion**

539    Identification and interpretation of genomic islands is challenging for two interrelated

540    reasons: combinations of evolutionary scenarios can lead to similar patterns of genomic

541    diversity and divergence, and widely used measures of divergence differ in sensitivity

542    and specificity to detect and differentiate between alternative scenarios. The original

543    explanation for island evolution, by which they arise in genomic regions where on-going

544    gene flow is effectively reduced due to divergent selection (Turner et al. 2005), has been

545    challenged and refined to include the various scenarios presented in Table 1. Whereas

546    empirical study will serve to resolve some of these issues, here we have used theoretical

547    approaches to explore the importance of linkage for establishment of new mutations

548    during an initial phase of local adaptation.

549

550    *Evolution of genomic islands by increased establishment probability*

551    There are three main factors that affect the likelihood of genomic islands evolving by

552    increased EP of linked mutations (mechanism 3.A): the relative increase in the EP of

553    linked vs. unlinked mutations, their absolute EP, and the size of the mutational target

554    within the linked region that experiences increased EP. The relative increase in EP for

555    linked vs. unlinked mutations will be greatest when there is very strong selection on a

556    single locus (the **B** locus) accompanied by high migration and weaker selection on other

557    loci (represented by the **A** locus here). If the distribution of fitness effects of new

558    mutations is gamma distributed, linked mutations will have much higher EP relative to

559    unlinked ones when the migration rate exceeds the average selection coefficient

560    substantially ($m >> \bar{a}$, Figure 3C). However, as the migration rate increases above $\bar{a}$, the

561    other factors lose importance: Both the absolute EP and the size of the region in which

562    linkage increases EP decrease rapidly (Figure 3A & D). Because the size of the

563    mutational target and the absolute EP both decrease, the waiting time for the

564    establishment of the first linked mutation increases in this region of parameter space. This

565    results in adaptation proceeding at a rate much slower than would occur in the absence of

566    migration (Figure 4). Thus, there is a 'goldilocks zone' where the balance between

567    migration and selection is just right for mechanism 3.A: migration is high enough that

568   unlinked mutations have low EP, but not so high that the waiting time for a linked

569   mutation to establish is very long. In other words, this zone is roughly between the

570   migration rate that maximizes the ratio of establishment probabilities in Figure 3C, and

571   the one that minimizes the time in Figure 4). In species with a continuous or stepping-

572   stone pattern of population structure spanning a broad environmental gradient, migration

573   rates may often be too geographically restricted to substantially constrain adaptation,

574   further limiting the potential importance of this mechanism.

575           In light of the sensitivity of this mechanism to the balance between migration and

576   selection, both demography and environment would have to be quite stable over

577   relatively long periods of time for evolution of genomic islands exclusively by increased

578   EP of linked mutations. In such cases, genomic islands would evolve most rapidly in

579   organisms with large population size, as this increases the total number of mutations

580   available for selection and therefore reduces the waiting time. Also, as the average

581   number of crossovers per chromosomes is found to be relatively insensitive to physical

582   chromosome length (Hillers and Villeneuve 2003), we would expect that organisms with

583   larger genomes have higher mutation rates per centimorgan, and that mechanism 3.A

584   might therefore be more important in such organisms. .

585           Two assumptions are implicit in the above discussion: the underlying mutational

586   target for a given polygenic trait is uniformly distributed throughout the genome and the

587   rate of recombination is homogeneous. If there are recombination coldspots around the B

588   locus due to segregating inversions or the previous fixation of recombination modifiers,

589   then increased EP due to linkage would extend over a greater physical distance and

590   therefore mutational target. Similarly, if functionally related loci tend to cluster together

591    in the genome (Nützmann & Osbourn 2014), then the rate of decay in the EP per base

592    pair would be unchanged, but the mutational target within this region would be increased.

593    Both of these non-homogeneous features of the genome can co-occur with the loci

594    involved in local adaptation by chance or as a result of previous bouts of local adaptation

595    that favoured the establishment or fixation of such modifiers or rearrangements

596    (mechanism 3.C). In either case, this could considerably increase the potential for the

597    evolution of genomic islands via increased probability of establishment of linked

598    mutations, as suggested by Yeaman (2013).

599          It is worth noting that Figure 4 provides a very rough comparison of the situation

600    for linked vs. unlinked mutations, as it uses the somewhat arbitrary $C_{95}$ to determine the

601    mutational target for linked mutations. This also assumes that the physical distance scales

602    linearly with the rate of recombination, which does not account for the reduction in

603    effective recombination rate when there are even numbers of crossovers between a pair of

604    loci in a given meiosis. Also, by simply contrasting the expected establishment times and

605    ignoring the variances, we do not provide a rigorous assessment of what would constitute

606    a statistically significant signature of an island (even a few clustered loci with many non-

607    clustered ones might be detectable as an island). Finally, our approach ignores the fact

608    that there might be multiple background loci that could independently initialise genomic

609    islands. The latter effect will depend on the DFE, and appropriate treatment of this would

610    require modelling $b$ as being drawn from that DFE, too. Comprehensively addressing

611    these issues is beyond the scope of this paper, but provides an obvious problem to address

612    in the future, and our conclusions throughout should be considered with these caveats in

613    mind.

614

*Evolution from new mutations vs. erosion and secondary contact*

616 Whereas mechanism 3.A depends explicitly on the mutation rate, the evolution of islands

617 via erosion following secondary contact (mechanism 3.D) is independent of mutation rate

618 and therefore can proceed much more rapidly (assuming that populations in allopatry

619 have had sufficient time to accumulate the differences that are eroded in secondary

620 contact). Here, we show that island formation via secondary contact and erosion of

621 divergence (mechanism 3.D) will be much more rapid than via establishment of new

622 mutations, even once a mutation that is destined to establish has already occurred (Figure

623 5). We also show that the expected time to island formation via secondary contact and

624 erosion scales approximately with the inverse of the effective migration rate ($1/m_e$; Figure

625 S4) for a wide range of population sizes. While we have not considered adaptation from

626 standing variation here, we suspect that the likelihood of island formation from standing

627 variation is similar to, if not lower than, that for new mutations for the same reasons as

628 we discussed above: in the parameter space where linkage is critical to establishment, the

629 absolute EP is low, so there would need to be many mutations present as standing

630 variation to give a signature. In addition, the fact that these mutations are segregating at

631 some intermediate, albeit potentially low, frequency implies that they have already

632 overcome the curse of stochastic loss. But this initial phase is exactly when linkage is

633 expected to convey the largest relative advantage. Mutations segregating as standing

634 variation may depend much less on this initial benefit. However, further work is

635 necessary to rigorously explore the contrast between standing and *de novo* variation in

636 this context, and special attention should be paid to the cause of standing variation, as

637    migration among populations inhabiting similar environments can introduce pre-adapted

638    variants.

639

640    *The absence of genomic islands*

641    While finding a strong signature of genomic islands can reveal much about the genomic

642    basis of adaptation, the lack of any observable peaks in $F_{ST}$ does not necessarily mean

643    that no adaptation is occurring. While the two-locus models covered here predict that no

644    local adaptation will occur when $m > m^*$, divergent adaptation at the phenotypic level can

645    be maintained at equilibrium through a quantitative genetic response mediated by very

646    small changes in allele frequency coupled with positive covariance among alleles,

647    provided sufficient genetic variance is maintained by mutation (Le Corre & Kremer

648    2003; Yeaman 2015). In such cases, divergence at the underlying loci may be transient,

649    and no signature of islands is expected, nor any significant peaks in $F_{ST}$. Alternatively, if

650    there are multiple alleles or haplotypes present at a single locus that yield the same

651    locally adapted phenotype, the signature of $F_{ST}$ would likely be greatly reduced,

652    especially at loci linked to the causal selected site(s). Studies of experimental evolution in

653    bacteria commonly find that many different amino acid residues within a given gene

654    evolve in response to the same selective pressure (e.g., Vogwill *et al.* 2014; Bailey *et al.*

655    2015). Similarly, a catalogue of the loci identified in adaptation has found that the same

656    locus often is involved repeatedly with different alleles (Martin & Orgogozo 2013). If

657    several such alleles at a single locus are maintained in a population, the expected

658    signatures of divergence would be weaker. Such patterns could likely evolve much more

659    readily with adaptation from standing variation. While we have restricted our study here

660      to two-deme, two-allele models, further study of more realistic population structures and

661      the possibility of haplotypes of equal fitness is necessary to more completely understand

662      their impact on statistical signatures of adaptation. Finally, it is worth noting that even in

663      cases where we would predict stable genomic islands to form based on our model, in

664      empirical studies these may not be statistically distinguishable from the genomic

665      background due to the highly stochastic nature of the genealogical process, mutation, and

666      recombination. This would be especially problematic in cases where migration between

667      populations is low and $F_{ST}$ at neutral sites is high (Feder & Nosil 2010). In such cases,

668      however, there may still be a genome-wide aggregate signal in the form of a negative

669      correlation between divergence and local recombination rates (Nachman & Payseur 2012;

670      Brandvain *et al.* 2014).

671

672      *Comparison to results of Yeaman (2013)*

673      The analysis presented here uses an approach similar to that of Yeaman (2013),

674      combining approximations for the EP with representations of the size of linked vs.

675      unlinked mutation targets, but extends this analysis and corrects a deficiency in

676      presentation. Because Yeaman (2013) numerically integrated the EP only over the

677      recombination rate but not the DFE, the analysis was restricted to considering the relative

678      advantage of linkage for a single mutation. As such, when migration was high enough

679      that the EP of an unlinked mutation became 0, the ratio of $\pi_L/\pi_U$ (referred to there as the

680      DHA, or 'Divergence Hitchhiking Advantage') became undefined. This was plotted as a

681      value of 0 in Yeaman's Figure 1, which is unfortunate as it suggests that the relative

682      advantage of linkage is low in this region of parameter space, when in fact it is infinitely

683    high and mutations of that size could not establish without linkage. Thus, inferences

684    based on Figure 1 from Yeaman (2013) discounted the contribution of potentially the

685    most important region of parameter space (i.e. the shaded region shown here in Figure

686    2A). Here, we correct this oversight by integrating across both recombination and the

687    DFE of possible mutations at the A locus, so that $\bar{\pi}_L$ always includes the contribution of

688    some mutations with non-zero EP (i.e., some mutations have $m^* > m$), and the mean EP,

689    $\bar{\pi}_L$ therefore never becomes 0 (Figure 2B). With this approach, we show how in the

690    region where $\bar{\pi}_L/\bar{\pi}_U$ tends to infinity, the advantage of linkage is counterbalanced by

691    reduced size of mutation target and increased time to establishment. Thus, our analysis

692    leads to the same qualitative conclusions as Yeaman (2013): that mechanism 3.A is

693    unlikely to be a broad explanation for the evolution of genomic islands independently of

694    some other factors affecting the distribution of the mutational target or the local rate of

695    recombination. These conclusions are insensitive to whether we use the splicing approach

696    (as in Yeaman 2013) or the two-type branching process (Aeschbacher & Bürger 2014) to

697    approximate the establishment probability. Although there are quantitative differences

698    between these two approximations if migration is relatively strong or recombination

699    relatively weak (Figure 1), there is good agreement between them in terms of the ratio of

700    establishment probabilities, $\bar{\pi}_L/\bar{\pi}_U$, and the 95% window size, $C_{95}$ (Figure S6).

701

702    *Conclusions and future directions*

703    Understanding the importance of the various mechanisms for the evolution of genomic

704    islands will inevitably require additional theory and more empirical work. Future

705    theoretical work could further explore the role of the DFE, study the effects of variable

706    selection strength at the background locus, and incorporate dominance, epistasis, and

707    other modes of selection (e.g. background selection). This would inform the development

708    of methods for robust inference about the mechanisms underlying genomic patterns of

709    diversity. In parallel, more comprehensive study of the extent of gene flow and

710    hybridization in natural populations is needed, as well as complementary lines of

711    evidence on gene function and phenotypic effect. This will inform the current debate

712    about the importance of gene flow in generating genomic islands (Cruickshank & Hahn

713    2014; Sætre 2014; Feulner *et al.* 2015; Monnahan *et al.* 2015) and help answer a number

714    of open questions: Is on-going gene flow commonly strong enough to make linkage

715    essential for divergence, or does gene flow only modulate evolving or previously existing

716    patterns? Are empirically observed islands produced by one or multiple sites under

717    divergent selection? How commonly do loci contribute to reproductive isolation

718    independent of adaptation? Do these islands reflect regions of reduced recombination as a

719    preadaptation to clustered architectures of traits under divergent selection, or did

720    clustered architectures and reduced local recombination rates evolve in response to such

721    selection? The distinction between preadaptation vs. response to selection may also be

722    seen as a question of timescale. For instance, in sticklebacks that repeatedly colonised

723    post-glacial lakes (Rogers *et al.* 2013), species may have encountered the same type of

724    environmental heterogeneity many times over deep evolutionary time. In such case, what

725    is now a pre-adaptation could be a consequence of previous response to selection for

726    reduced recombination between locally adapting loci. From the theoretical arguments

727    outlined here, it seems unlikely that increased establishment probability due to linkage

728    will provide a broad explanation for the evolution of genomic islands without some sort

729  of preadaptation, either in local recombination rate or the genomic distribution of

730  mutational targets. Comparative genomic studies may provide the best means to identify

731  how commonly such genome-scale changes have shaped the chromosomal landscape

732  upon which local adaptation and reproductive isolation are built.

733
734
735  **Acknowledgements**

744
745  **References**

746  Aeschbacher S, Bürger R (2014) The effect of linkage on establishment and survival of
747      locally beneficial mutations. *Genetics*, **197**, 317–336.

748  Akerman A, Bürger R (2014) *The consequences of gene flow for local adaptation and*
749      *differentiation: a two-locus two-deme model*.

750  Bailey SF, Rodrigue N, Kassen R (2015) The Effect of Selection Environment on the
751      Probability of Parallel Evolution. *Molecular Biology and Evolution*, **32**, 1436–1448.

752  Bank C, Bürger R, Hermisson J (2012) The limits to parapatric speciation: Dobzhansky-
753      Muller incompatibilities in a continent-islands model. *Genetics*, **191**, 845–863.

754  Barton N, Bengtsson BO (1986) The barrier to genetic exchange between hybridising
755      populations. *Heredity*, **57 ( Pt 3)**, 357–376.

756  Barton N, Hewitt GM (1989) Adaptation, speciation and hybrid zones. *Nature*, **341**, 497–
757      503.

758  Brandvain Y, Kenney AM, Flagel L, Coop G, Sweigart AL (2014) Speciation and
759      Introgression between Mimulus nasutus and Mimulus guttatus. *PLoS Genetics*, **10**,
760      e1004410.

761  Bürger R, Akerman A (2011) The effects of linkage and gene flow on local adaptation: A
762      two-locus continent-island model. *Theoretical Population Biology*, **80**, 272–288.

763  Charlesworth B (1979) Selection on recombination in. , 581–589.

764  Charlesworth B (1998) Measures of divergence between populations and the effect of
765      forces that reduce variability. *Molecular biology and evolution*, **15**, 538–543.

766 Charlesworth B, Morgan MT, Charlesworth D (1993) The effect of deleterious mutations
767     on neutral molecular variation. *Genetics*, **134**, 1289–1303.

768 Charlesworth B, Nordborg M, Charlesworth D (1997) The effects of local selection,
769     balanced polymorphism and background selection on equilibrium patterns of genetic
770     diversity in subdivided populations. *Genetical research*, **70**, 155–174.

771 Le Corre V, Kremer A (2003) Genetic variability at neutral markers, quantitative trait loci
772     and trait in a subdivided population under selection. *Genetics*, **164**, 1205–1219.

773 Cruickshank TE, Hahn MW (2014) Reanalysis suggests that genomic islands of
774     speciation are due to reduced diversity, not reduced gene flow. *Molecular Ecology*,
775     **23**, 3133–3157.

776 Excoffier L, Hofer T, Foll M (2009) Detecting loci under selection in a hierarchically
777     structured population. *Heredity*, **103**, 285–298.

778 Feder JL, Gejji R, Yeaman S, Nosil P (2012) Establishment of new mutations under
779     divergence and genome hitchhiking. *Philosophical Transactions of the Royal
780     Society B: Biological Sciences*, **367**, 461–474.

781 Feder JL, Nosil P (2010) The efficacy of divergence hitchhiking in generating genomic
782     islands during ecological speciation. Evolution, **64**, 1729–1748.

783 Feulner PGD, Chain FJJ, Panchal M *et al.* (2015) Genomics of divergence along a
784     continuum of parapatric population differentiation. *PLoS genetics*, **11**, e1004966.

785 Fisher R (1930) *The Genetical Theory of Natural Selection*. Clarendon, Oxford.

786 Gillespie JH (2000) The neutral theory in an infinite population. *Gene*, **261**, 11–18.

787 Griffiths RC (1981) The number of heterozygous loci between two randomly chosen
788     completely linked sequences of loci in two subdivided population models. *Journal
789     of Mathematical Biology*, **12**, 251–261.

790 Haldane JBS (1930) A mathematical theory of natural and artificial selection. VI.
791     Isolation. *Mathematical Proceedings of the Cambridge Philosophical Society*, **26**,
792     220–230.

793 Harr B (2006) Genomic islands of differentiation between house mouse subspecies. ,
794     730–737.

795 Hermisson J (2009) Who believes in whole-genome scans for selection? *Heredity*, **103**,
796     283–284.

797 Hillers KJ, Villeneuve AM (2003) Chromosome-wide control of meiotic crossing over in
798     C. elegans. *Current Biology*, **13**, 1641–1647.

799 Kirkpatrick M (2006) Chromosome Inversions, Local Adaptation and Speciation.
800     *Genetics*, **173**, 419–434.

801 Lenormand T, Otto SP (2000) The evolution of recombination in a heterogeneous
802     environment. *Genetics*, **156**, 423–438.

803 Lindtke D, Buerkle CA (2015) The genetic architecture of hybrid incompatibilities and

804      their effect on barriers to introgression in secondary contact. *Evolution*, **69**, 1987-
805      2004.

806 López E, Pradillo M, Oliver C *et al.* (2012) Looking for natural variation in chiasma
807      frequency in Arabidopsis thaliana. *Journal of Experimental Botany*, **63**, 887–894.

808 Martin A, Orgogozo V (2013) The loci of repeated evolution: a catalog of genetic
809      hotspots of phenotypic adaptation. *Evolution*, **67**, 1235-1250.

810 Maynard Smith J, Haigh J (1974) The hitch-hiking effect of a favourable gene. *Genetical*
811      *Research Cambridge*, **23**, 23–35.

812 Monnahan PJ, Colicchio J, Kelly JK (2015) A genomic selection component analysis
813      characterizes migration-selection balance. *Evolution*, n/a–n/a.

814 Nei M (1973) Analysis of gene diversity in subdivided populations. *Proceedings of the*
815      *National Academy of Sciences of the United States of America*, **70**, 3321–3323.

816 Nei M, Li WH (1979) Mathematical model for studying genetic variation in terms of
817      restriction endonucleases. *Proceedings of the National Academy of Sciences of the*
818      *United States of America*, **76**, 5269–5273.

819 Nachman, MW. and Payseur, BA. (2012) Recombination rate variation and speciation:
820 theoretical predictions and empirical results from rabbits and mice. *Phil. Trans. R. Soc. B*
821 367:409–421.
822
823 Noor MAF, Grams KL, Bertucci LA, Reiland J (2001) Chromosomal inversions and the
824      reproductive isolation of species. , **2001**.

825 Nordborg M (1997) Structured coalescent processes on different time scales. *Genetics*,
826      **146**, 1501–1514.

827 Nosil P, Funk DJ, Ortiz-Barrientos D (2009) Divergent selection and heterogeneous
828      genomic divergence. *Molecular Ecology*, **18**, 375–402.

829 Notahara M (1990) The coalescent and the genealogical process in geographically
830      structured population. *Journal of Mathematical Biology*, **29**, 59–75.

831 Nützmann H-W, Osbourn A (2014) Gene clustering in plant specialized metabolism.
832      *Current Opinion in Biotechnology*, **26**, 91–99.

833 Petry D (1983) The Effect on Neutral Gene Flow of Selection at a Linked Locus. , **313**,
834      300–313.

835 Renaut S, Grassa CJ, Yeaman S *et al.* (2013) Genomic islands of divergence are not
836      affected by geography of speciation in sunflowers. *Nature Communications*, **4**,
837      1827.

838 Rieseberg LH (2001) Chromosomal rearrangements and speciation. *Trends in Ecology*
839      *and Evolution*, **16**, 351–358.

840 Rogers SM, Bowles E, Mee JA The consequences of genomic architecture on ecological
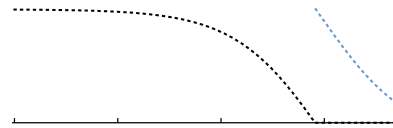841      speciation in postglacial fishes. , 1–25.

842   Sætre G-P (2014) Genome scans and elusive candidate genes: detecting the variation that
843       matters for speciation. *Molecular ecology*, **23**, 4677–8.

844   Slatkin M (1991) Inbreeding coefficients and coalescence times. *Genetical Research*
845       *Cambridge*, **58**, 167–175.

846   Slatkin M, Excoffier L (2012) Serial founder effects during range expansion: A spatial
847       analog of genetic drift. *Genetics*, **191**, 171–181.

848   Strasburg JL, Sherman NA, Wright KM *et al.* (2012) What can patterns of differentiation
849       across plant genomes tell us about adaptation and speciation ? What can patterns of
850       differentiation across plant genomes tell us about adaptation and speciation ?

851   Takahata N (1988) The coalescent in two partially isolated diffusion populations.
852       *Genetical Research*, **52**, 213–222.

853   Trans P, Lond RS, Barton NH (2000) Genetic hitchhiking Genetic hitchhiking. , 1553–
854       1562.

855   Turner TL, Hahn MW, Nuzhdin S V (2005) Genomic Islands of Speciation in Anopheles
856       gambiae. , **3**.

857   Via S, West J (2008) The genetic mosaic suggests a new role for hitchhiking in
858       ecological speciation. *Molecular Ecology*, **17**, 4334–4345.

859   Vogwill T, Kojadinovic M, Furió V, MacLean RC (2014) Testing the Role of Genetic
860       Background in Parallel Evolution Using the Comparative Experimental Evolution of
861       Antibiotic Resistance. *Molecular biology and evolution*, **31**, 3314–3323.

862   Wilkinson-Herbots H (1998) Genealogy and subpopulation differentiation under various
863       models of population structure. *Journal of Mathematical Biology*, **37**, 535–585.

864   Wilkinson-Herbots HM (2008) The distribution of the coalescence time and the number
865       of pairwise nucleotide differences in the "isolation with migration" model.
866       *Theoretical Population Biology*, **73**, 277–288.

867   Wilkinson-Herbots HM (2012) The distribution of the coalescence time and the number
868       of pairwise nucleotide differences in a model of population divergence or speciation
869       with an initial period of gene flow. *Theoretical Population Biology*, **82**, 92–108.

870   Wright S (1931) Evolution in Mendelian populations. *Genetics*, **16**, 97–159.

871   Wright S (1943) Isolation by distance. *Genetics*, **28**.

872   Yeaman S (2013) Genomic rearrangements and the evolution of clusters of locally
873       adaptive loci. *Proceedings of the National Academy of Sciences of the United States*
874       *of America*, **110**, E1743–51.

875   Yeaman S (2015) (in press) Local adaptation by small-effect alleles. *Am Nat*, **186**.

876   Yeaman S, Otto SP (2011) Establishment and maintenance of adaptive genetic
877       divergence under migration, selection, and drift. *Evolution*, **65**, 2123–2129.

878   Yeaman S, Whitlock MC (2011) The Genetic Architecture of Adaptation under
879       Migration-Selection Balance. *Evolution*, **65**, 1897–1911.
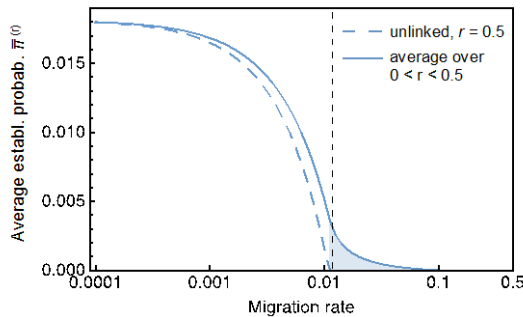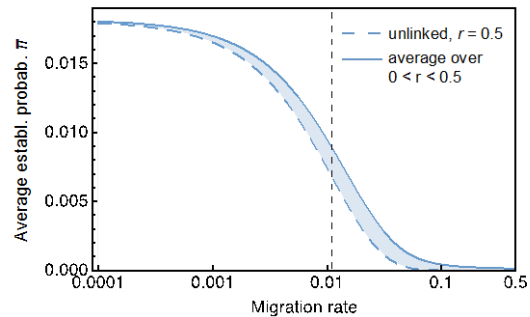
**Figures**

881
882

A

883
884    **Figure 1**. Comparison of three approximations to the establishment probability of a new
885    linked mutation. The establishment probability based on the two-type branching process
886    (solid lines), a slightly-supercritical two-type branching process assuming $a, m \ll b, r$
887    (dashed lines), and the splicing approach (dotted lines) as a function of migration rate (A)
888    and recombination rate (B). A) $r = 0.01$; B) $m = 0.01$. In both cases, $b = 0.1$
889
890
891

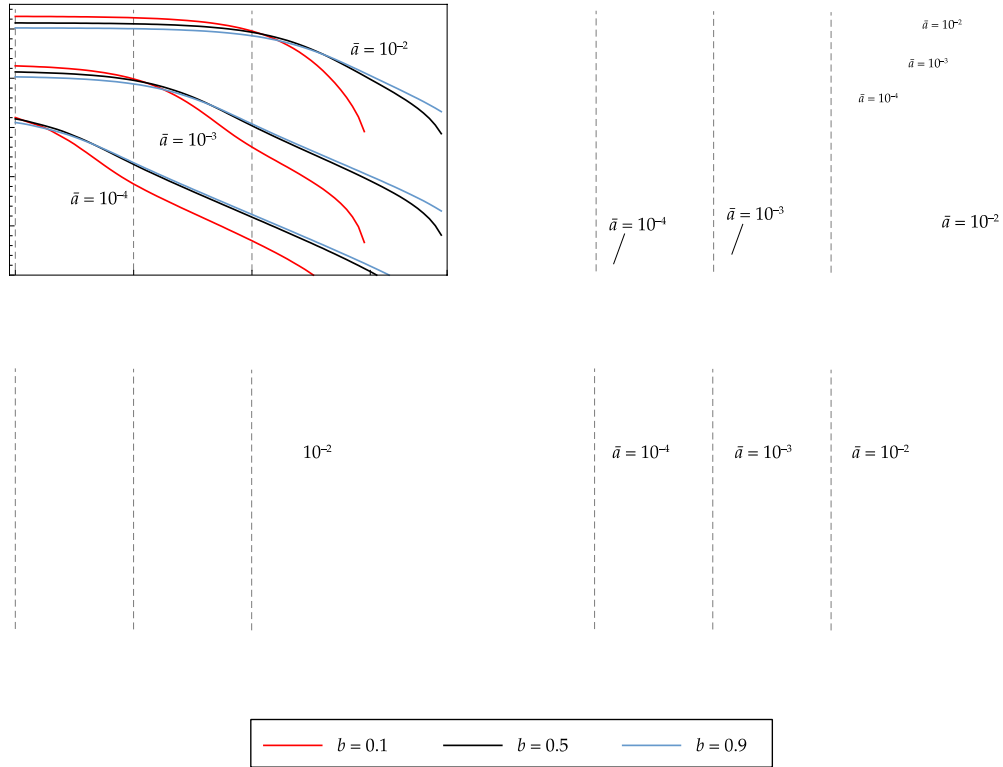A                                                    B



892
893
894    **Figure 2.** The effect of linkage on the average establishment probability of locally
895    beneficial mutations. A) Contrast between an unlinked ($\bar{\pi}_U^{(r)}$, dashed curve) and an
896    average linked ($\bar{\pi}_L^{(r)}$, solid curve) mutation of effect $a = 0.01$. B) As in A), but after
897    averaging over an exponential distribution of fitness effects (DFE) with mean $\bar{a} = 0.01$
898    (dashed for $\bar{\pi}_U$ vs. solid for $\bar{\pi}_L$). In both panels, the selection coefficient at the
899    background locus is $b = 0.1$, and the vertical dashed line indicates the critical migration
900    rate below which a single unlinked mutation of effect $a = 0.01$ can be established.
901    Results are shown for the two-type branching process.
902

$\bar{a} = 10^{-2}$

$\bar{a} = 10^{-3}$

$\bar{a} = 10^{-4}$

$\bar{a} = 10^{-4}$    $\bar{a} = 10^{-3}$    $\bar{a} = 10^{-2}$

$\bar{a} = 10^{-2}$

$\bar{a} = 10^{-3}$

$\bar{a} = 10^{-4}$

$\bar{a} = 10^{-2}$

$10^{-2}$    $\bar{a} = 10^{-4}$    $\bar{a} = 10^{-3}$    $\bar{a} = 10^{-2}$
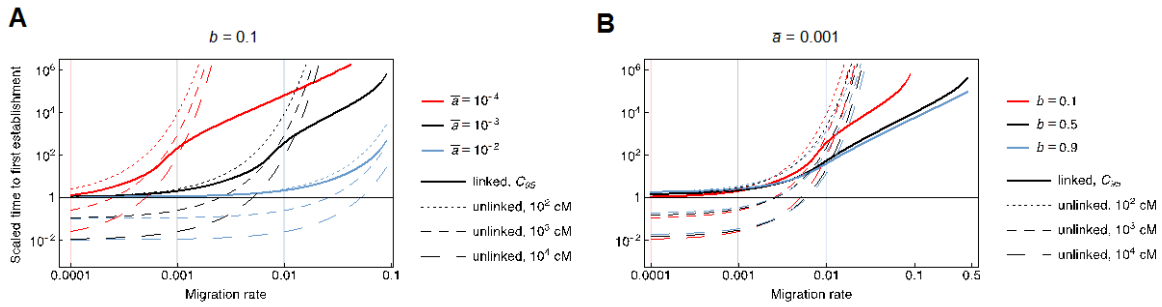
— $b = 0.1$    — $b = 0.5$    — $b = 0.9$

904
905
906 **Figure 3.** The effect of the migration rate on the establishment probability and window
907 size. The average establishment probability $\bar{\pi}_L$ as a function of the migration rate for
908 various strengths of selection on the $\log_{10}$ (A) and natural scale (B), and relative to the
909 average establishment probability of an unlinked mutation, $\bar{\pi}_L/\bar{\pi}_U$ (C). D) The size of the
910 window within which 95% of all successfully establishing linked *de novo* mutations
911 occur ($C_{95}$). All panels show results for the two-type branching process with different
912 values of $\bar{a}$ and $b$, and $k = 1$ (exponential DFE).
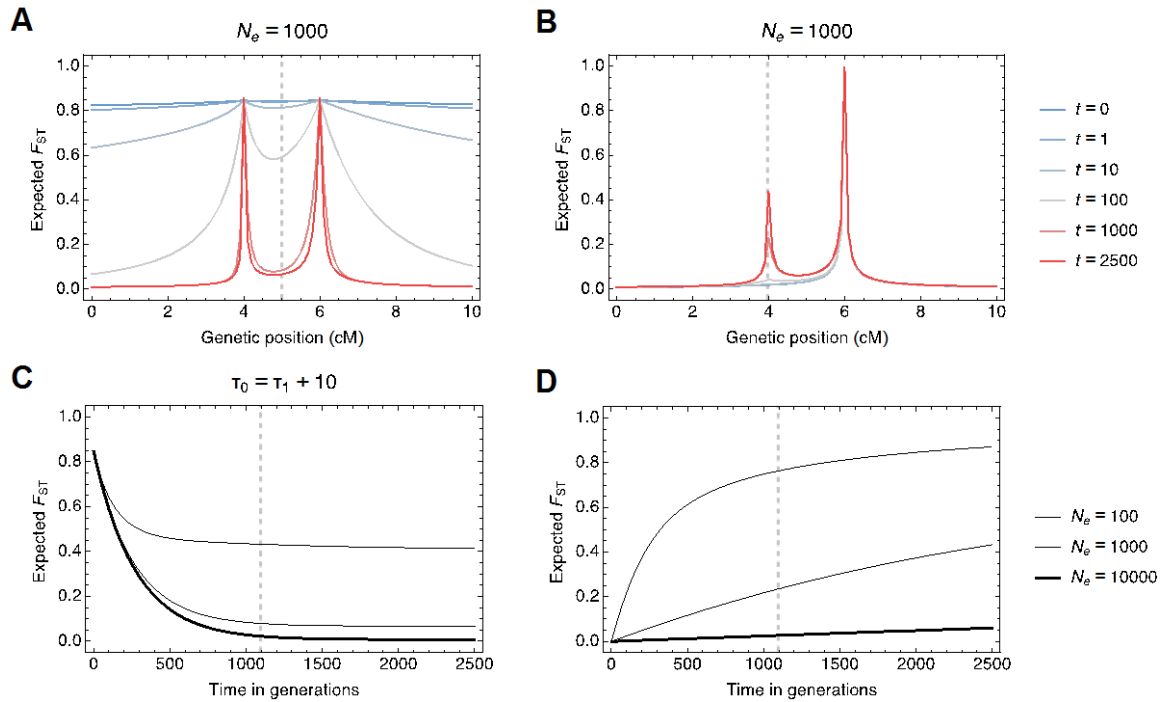913
914

915



916
917
918 **Figure 4.** Waiting times to the first establishing linked and unlinked locally beneficial
919 mutation. Expected waiting times are shown as a function of the migration rate for linked
920 (solid curves) and unlinked (dashed curves) mutations for various mean selection
921 coefficients $\bar{a}$ at the focal allele (A) and different selection coefficients $b$ at the
922 background locus (B). Times for unlinked mutations are shown assuming three different
923 values of the unlinked mutational target $C_U$ (different degrees of dashing for 100, 1000,
924 and 10000 cM). Waiting times are scaled by the corresponding substitution rate in a
925 completely isolated panmictic population of the same size (see text for details). In both
926 panels, $k = 1$ (exponential DFE; see Figure S3 for other shapes), and predictions are
927 based on the two-type branching process. Light grey lines included to facilitate
928 comparison between parameter combinations.
929
930

931
932
933
934 **Figure 5.** Dynamics of $F_{ST}$ upon secondary contact and during the rise of a new selected
935 allele in finite populations. A, B) Divergence between two demes of size $N_e = 1000$ in
936 terms of $F_{ST}$ at a neutral site as a function of its genetic position at various points in time.
937 A) Formation of a two-peak island by erosion of neutral divergence in the neighbourhood
938 of two selected loci A and B positioned at 4 and 6 cM. Before secondary contact, A and B
939 have undergone adaptive divergence over 20 $N_e$ generations in allopatry. B) Formation of
940 a two-peak island through the rise of a locally beneficial *de novo* mutation at locus A in
941 the face of gene flow and in the vicinity of a previously established migration–selection
942 polymorphism at B. C, D) Dynamics of $F_{ST}$ at a neutral site in the centre of the two-peak
943 island (C) and at the position of the rising selected locus (D) for different effective
944 population sizes $N_e$. Vertical lines indicate the inverse of the effective rate of gene flow,
945 $1/m_e$, in the centre of the island, as a deterministic approximation to the waiting time for
946 the erosion of neutral divergence. In all panels, $F_{ST}$ is approximated as a ratio of
947 coalescence times under the appropriate demographic model, with actual migration rates
948 replaced by the deterministic approximation to the effective rate of gene flow, $m_e$, at the
949 position of the neutral site (see SI for details). Parameters are $a = 0.01$, $b = 0.1$, and $m =$
950 0.02.